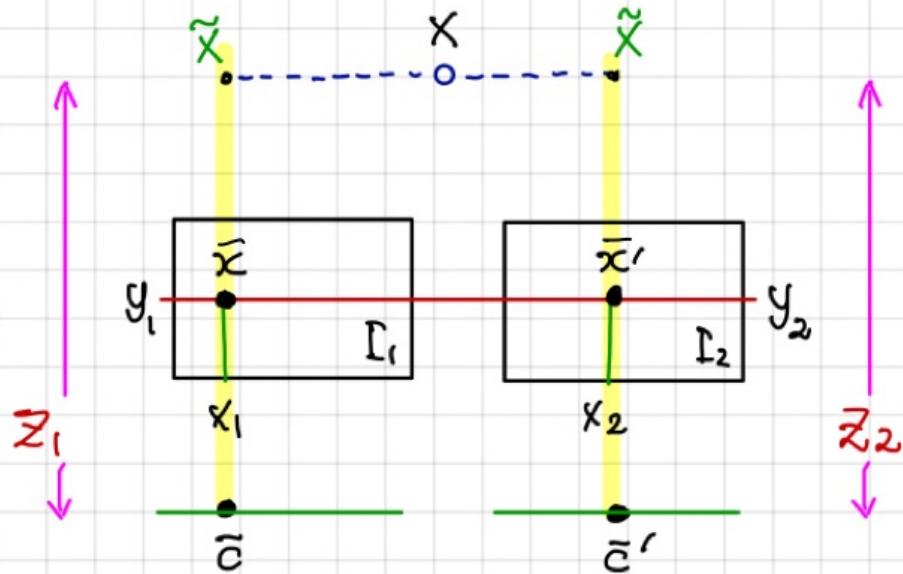
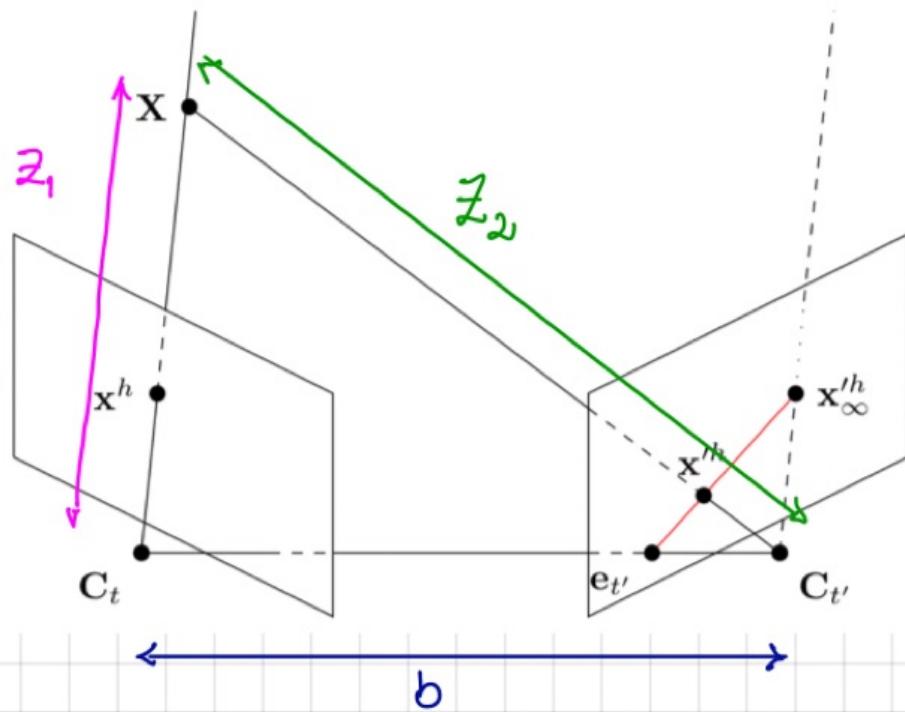


# LECTURE 8 : DEPTH FROM VIDEO

For a pair of rectified images , there is only one disparity or depth map :



Though, before being rectified, there are 2 depth maps :



Q: How to compute disparity of non-rectified images?

A: Two ways :

- (1) Rectify the images  $\rightarrow$  Epipolar rectification
- (2) Using epipolar geometry



Epipolar-geometry implicit rectification is preferable,  
since less steps are required in the computation.



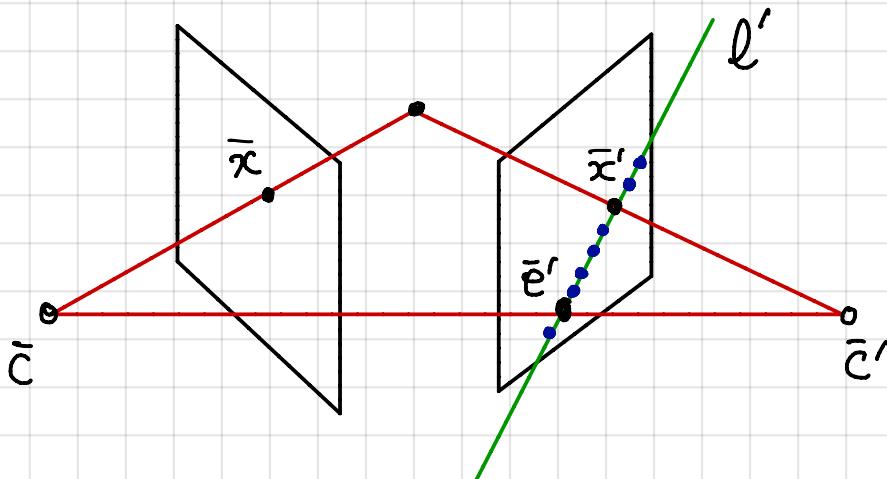
The basic formulation:

$$d = \frac{1}{z} \quad ; \quad \begin{array}{l} d = \text{disparity} \\ z = \text{depth} \end{array}$$

$$\bar{x}' \sim K' R' R^T K^{-1} \bar{x} + d K' R' (\bar{c} - \bar{c}')$$

Given  $\bar{x}$ ,  $P = K[R|E]$  and  $P' = K'[R'|E']$ , we want to  
find  $d$  using the formula, such that  $I(\bar{x}) = I'(\bar{x}')$ .

The formula can generate all points on the green line below,  
when we change the value of  $d$ :

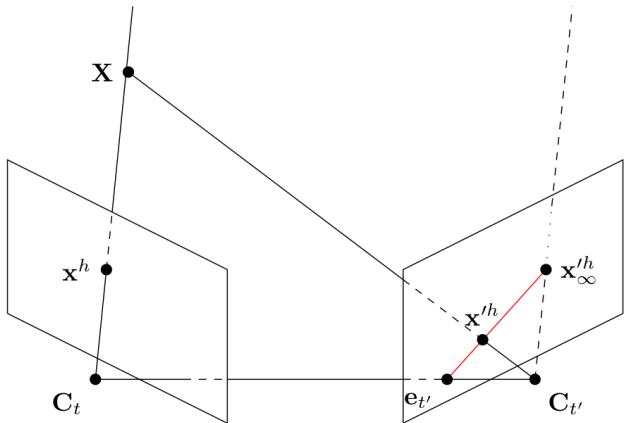


## [4] Derivation of Epipolar Rectification

To understand how we can get this equation:

$$\bar{x}' \sim K' R' R^T K^{-1} \bar{x} + d K' R' (\bar{c} - \bar{c}'),$$

consider:



$$\bar{x}^h = P X = K[R|\bar{t}] X$$

Assume the existence of  $X_\infty$  along the line that passes through  $\bar{x}$ , where  $X_\infty$  is located at infinity:

$$X_\infty = \begin{bmatrix} x \\ y \\ z \\ 0 \end{bmatrix} = \begin{bmatrix} \hat{x}_\infty \\ 0 \end{bmatrix}$$

$$\underset{3 \times 1}{\bar{x}^h} = \underset{3 \times 3}{K[R|\bar{t}]} \underset{3 \times 3}{X_\infty} = \underset{3 \times 3}{K[R|\bar{t}]} \underset{\frac{3 \times 4}{3 \times 1}}{\begin{bmatrix} \hat{x}_\infty \\ 0 \end{bmatrix}} = \underset{3 \times 3}{K R \hat{x}_\infty} \underset{3 \times 1}{\hat{x}_\infty}$$

Thus:

$$\hat{x}_\infty = R^T K^{-1} \bar{x}^h$$

Projecting  $X_\infty$  onto the right image:

$$\begin{aligned} \bar{x}'^h &= P' X_\infty = K' [R' | \bar{t}'] \begin{bmatrix} \hat{x}_\infty \\ 0 \end{bmatrix} = K' R' \hat{x}_\infty \\ &= K' R' R^T K^{-1} \bar{x}^h \end{aligned}$$

What we want to find is  $\bar{x}^h$ , and not  $\bar{x}'^h$ .

Basic idea of finding  $\bar{x}^h$ :

To find  $\bar{x}^h$ , we can start the search from  $\bar{x}_\infty^h$ . Then from  $\bar{x}_\infty^h$ , we can trace along the epipolar line in the direction of the epipolar point,  $\bar{e}_{t'}$ :

$$\begin{aligned}\bar{e}_{t'} &= P' \bar{c}_t = K' [R' | \bar{E}'] \bar{c}_t \\ &= K' R' [\underbrace{I}_{3 \times 3} \underbrace{[R]^{-1} \bar{E}'}_{3 \times 1}] \bar{c}_t \\ &= K' R' [\bar{c}_t - \bar{c}_{t'}]\end{aligned}$$

Recall:

$$\bar{c}_{t'} = -(R')^{-1} \bar{E}'$$

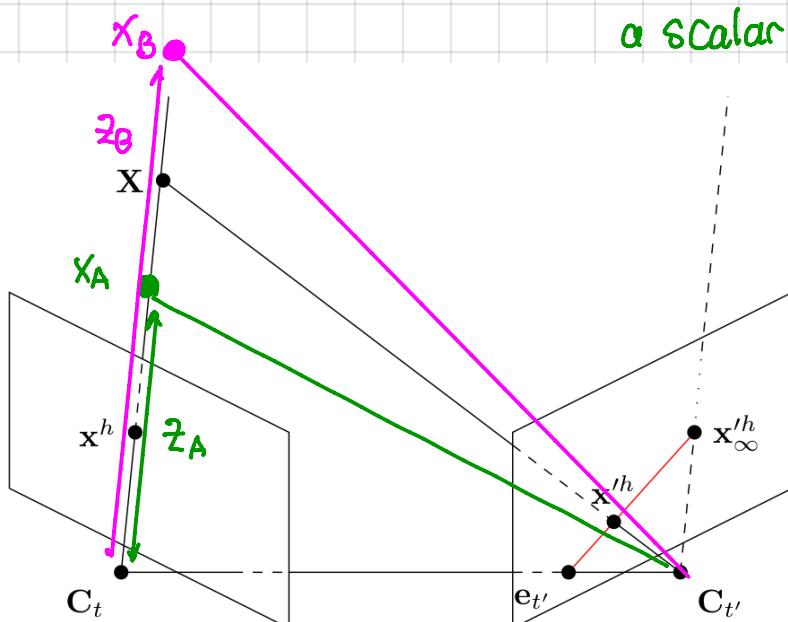
Therefore:

$$\bar{x}^h \underset{3 \times 1}{\sim} \bar{x}_\infty^h + d \bar{e}_{t'} \rightarrow \text{the line parametric equation}$$

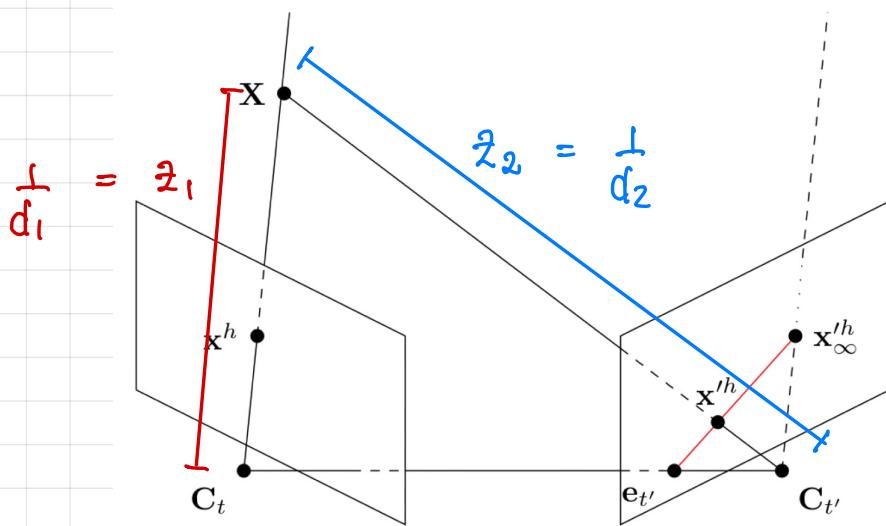
$$\bar{x}^h \sim K' R' R^T K^{-1} \bar{x}^h + d \bar{e}_{t'}$$

$$\bar{x}^h \underset{3 \times 1}{\sim} \underset{3 \times 3}{K' R'} \underset{3 \times 3}{R^T} \underset{3 \times 3}{K^{-1}} \underset{3 \times 1}{\bar{x}^h} + d \underset{3 \times 3}{K' R'} [\bar{c}_t - \bar{c}_{t'}]$$

a scalar value



The shorter the depth ( $z_A$ ), the larger the value of  $d$ , vice versa. This also indicates  $d$  determines the depth of point  $x$ .



$d$  = disparity  
 $z$  = depth

In the figure, with respect to  $X$ , there are 2 depths :

- $z_1 = 1/d_1$  is the depth from  $C_t$  to  $X$
- $z_2 = 1/d_2$  is the depth from  $C_{t'}$  to  $X$

$$\bar{x}^h \sim \bar{x}'^h + d_1 \bar{e}_{t'} \quad \rightarrow \text{You need to normalize } x^h \\ \begin{matrix} 3 \times 1 & 3 \times 1 & 3 \times 1 \end{matrix} \quad \text{so that: } \begin{pmatrix} x \\ y \\ z \end{pmatrix} \rightarrow \begin{pmatrix} x/z \\ y/z \\ 1 \end{pmatrix}$$

Q: How to justify that this equation is correct ?

A: Consider the following cases :

If  $d=0$ , meaning  $z=\infty$  :

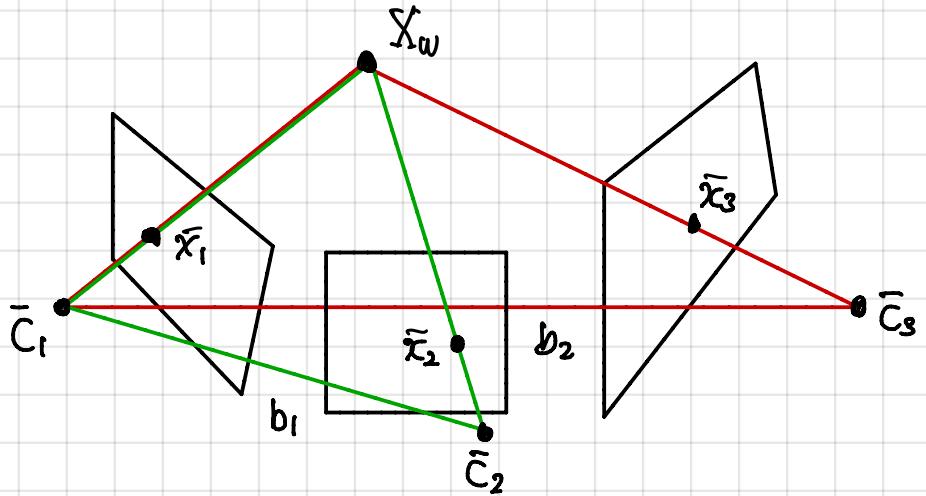
$$\bar{x}^h = \bar{x}'^h$$

If  $d=\infty$ , meaning  $z=0$  (the 3D point is at  $\bar{C}_t$ ) :

$$\begin{aligned} \bar{x}^h &= \bar{x}'^h + \infty \bar{e}' \\ &= \bar{e}' \end{aligned} \quad \left. \right\} \text{see the illustration above}$$

If  $0 < d < \infty$ , then  $\bar{x}^h$  must lie on the epipolar line !

## [5] Epipolar Rectification for Multiple Images



[0] For image 1 & 2 :

$$\bar{x}_2 = \bar{x}_{200} + d_{12} \quad \bar{e}_2 = K_2 R_2 R_1^T K_1^{-1} \bar{x}_1 + d_{12} K_2 R_2 (\bar{c}_1 - \bar{c}_2)$$

[0] For image 1 & 3 :

$$\bar{x}_3 = \bar{x}_{300} + d_{13} \quad \bar{e}_3 = K_3 R_3 R_1^T K_1^{-1} \bar{x}_1 + d_{13} K_3 R_3 (\bar{c}_1 - \bar{c}_3)$$



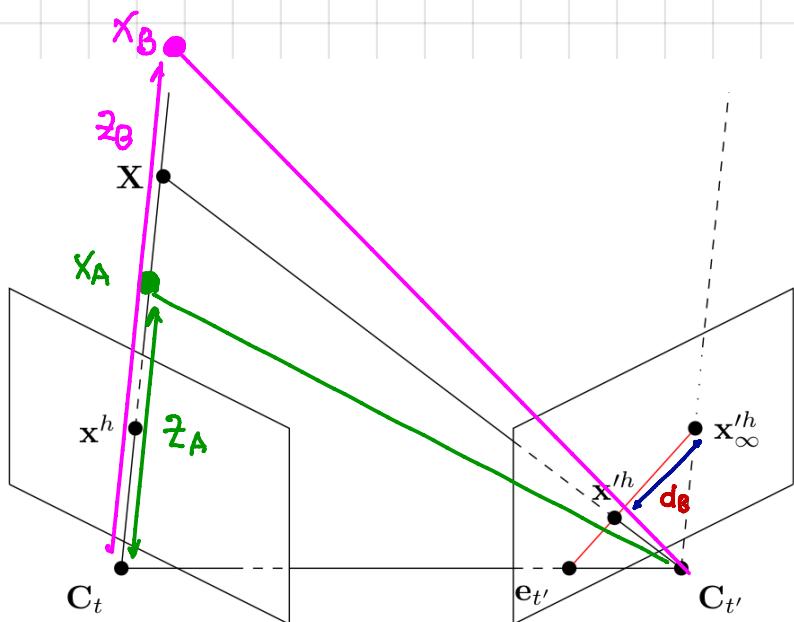
$d_{12} = d_{13}$

Implying :  $I_2(\bar{x}_2) = I_1(\bar{x}_1)$

$I_3(\bar{x}_3) = I_1(\bar{x}_1)$

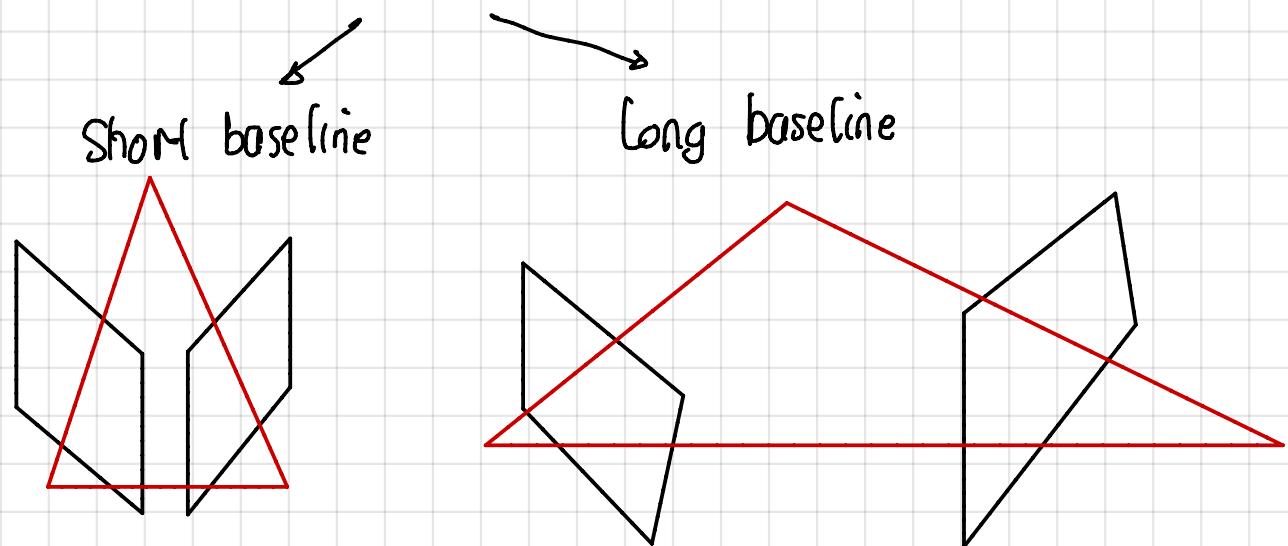
Q : Why  $d_{12} = d_{13}$  ?

A : Because  $d_{12}$  &  $d_{13}$  indicate the same depth of  $X$  from  $C_1$ .



## [6] Multiple Baseline Stereo = Depth from Video

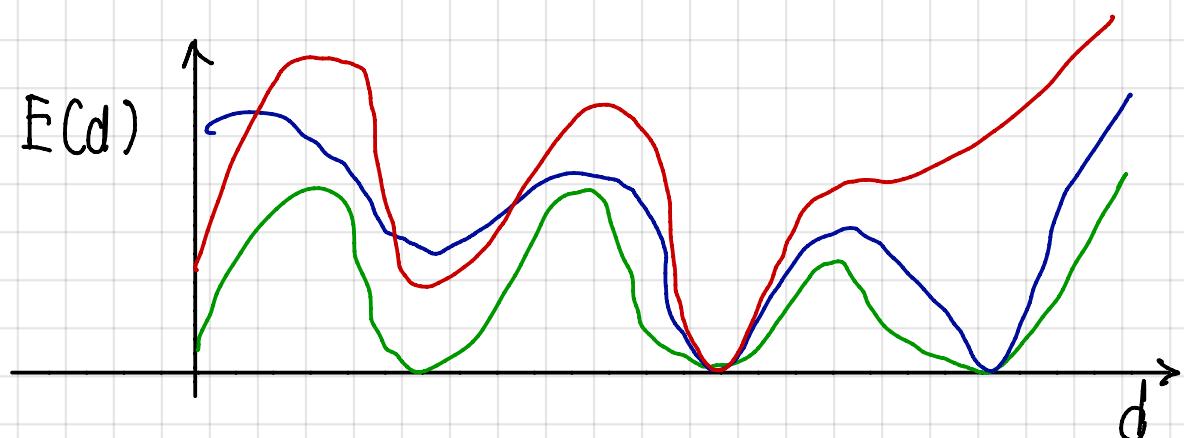
Disparity values depends on the baseline



For estimating disparity, a longer baseline is generally better, since  $b$  in the equation acts like a magnifier. However, longer baselines tend to suffer from occlusions.



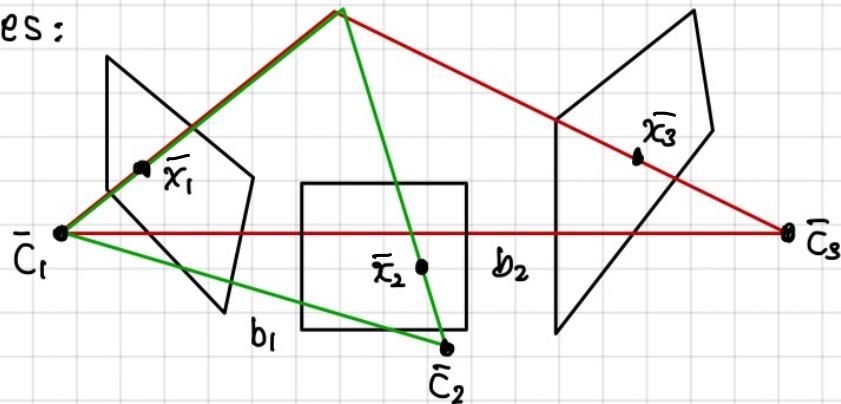
Multiple baselines allow us to get the benefits of short and long baselines:



The green line indicates that there are many candidates of  $d$  that make  $E(d)=0 \Rightarrow$  ambiguity!

## [7] Depth from Video (Multiple Baselines)

Multiple baselines:



Two main steps

(1) Disparity Initialization

(2) Bundle Optimization

Purpose: To have an initial depth map of each frame.

Input:  $\{I_t, P_t\}_{t=1}^N$ ;  $t = 1 \dots N$

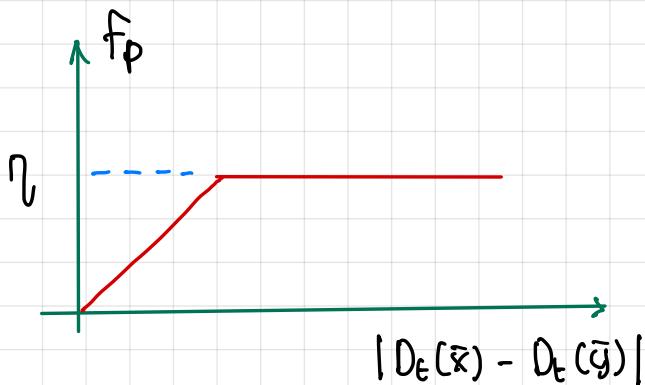
Output:  $\{D_t\}_{t=1}^N$  (the disparity maps)

$N$  = the number frames.

MAF:

$$D_t^{\text{init}} = \underset{\{d_{\min} \dots d_{\max}\}}{\operatorname{argmin}} \sum_x \left[ l - U(\bar{x}) f_d^{\text{init}}(\bar{x}, D_t(\bar{x})) + \sum_{\bar{y} \in N\bar{x}} \lambda(\bar{x}, \bar{y}) f_p(D_t(\bar{x}), D_t(\bar{y})) \right]$$

• Prior term:  $f_p(D_t(\bar{x}), D_t(\bar{y})) = \min(|D_t(\bar{x}) - D_t(\bar{y})|, \eta)$



To prevent  $f_p$  from being too large, since the difference between two disparity values shouldn't be that large.

Main reference: "Consistent depth map recovery from a video sequence."  
TPAMI, 2009.

MRF:

$$D_t^{\text{init}} = \underset{\{d_{\min} \dots d_{\max}\}}{\operatorname{argmin}} \sum_{\bar{x}} \left[ l - U(\bar{x}) f_d^{\text{init}}(\bar{x}, D_t(\bar{x})) + \sum_{\bar{y} \in N\bar{x}} f_p(D_t(\bar{x}), D_t(\bar{y})) \right]$$

- Data term:

Use all other frames (prev. & subsequent frames)

$$f_d^{\text{init}}(\bar{x}, D_t(\bar{x})) = \sum_{t'}^N f_c(\bar{x}, D_t(\bar{x}), I_t, I_{t'})$$

e.g.:  $t = 1$   
 $t' = 2, 3, \dots, N$

where:

$$f_c(\bar{x}, d, I_t, I_{t'}) = \frac{\sigma_c}{\sigma_c + |I_t(\bar{x}) - I_{t'}(l_{t,t'}(\bar{x}, d))|}$$

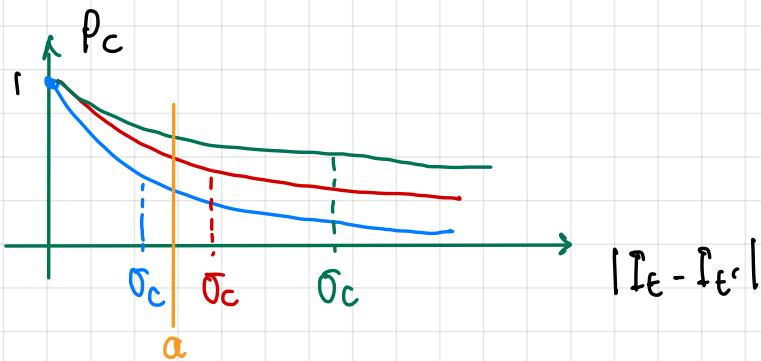
photo consistency

$l_{t,t'}(\bar{x}_t, d)$  means:

$$\bar{x}_{t'} = K_{t'} R_{t'} R_t^T K_t^{-1} \bar{x}_t + d K_{t'} R_{t'} [\bar{c}_t - \bar{c}_{t'}]$$

$3 \times 1 \quad 3 \times 3 \quad 3 \times 3 \quad 3 \times 3 \quad 3 \times 1 \quad 1 \times 1 \quad 3 \times 3 \quad 3 \times 3 \quad 3 \times 1$

$\sigma_c$  implies a variable to control the tolerance of the photo consistency:



If  $|I_t - I_{t'}| = \alpha$ , then different values of  $\sigma_c$  will produce different  $P_c$ .

The smaller  $\sigma_c$  (blue line) penalize  $|I_t - I_{t'}|$  more. Meaning, even though  $|I_t - I_{t'}|$  is relatively small,  $P_c$  will be smaller compared with the other values of  $\sigma_c$  (red and green lines).

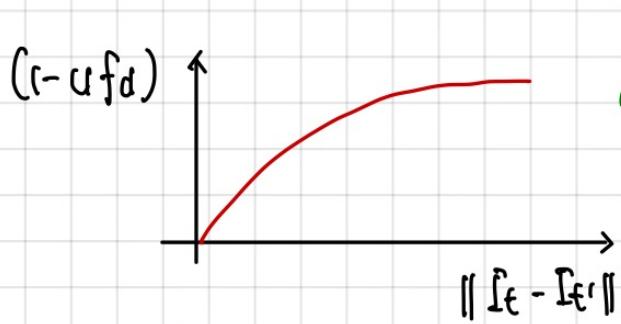
MRF :

$$D_t^{\text{init}} = \underset{\{d_{\min} \dots d_{\max}\}}{\operatorname{argmin}} \sum_{\bar{x}} \left[ (1 - U(\bar{x})) f_d^{\text{init}}(\bar{x}, D_E(\bar{x})) + \sum_{\bar{y} \in N_{\bar{x}}} \lambda(x, y) f_p(D_t(\bar{x}), D_t(\bar{y})) \right]$$

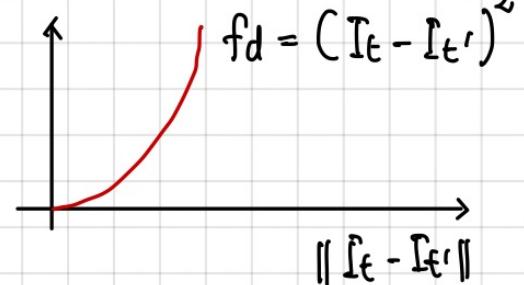
$$U(\bar{x}) = \frac{1}{\max_{D_E(\bar{x})} f_d(\bar{x}, D(\bar{x}))}$$

→ the normalization factor.  
To ensure the values ranging from 0 to 1.

The graph of the overall data term:



more robust than  
our prev. cost  
function



If there is noise  $\|I_t - I_{t'}\|$  will be large, hence  $(I_t - I_{t'})^2$  will be also large, though the actual cost shouldn't be that large.

Using  $(1 - U f_d)$  will reduce the effect of noise in the optimization.

- Weighting factor  $\lambda$ :

MRF:

$$D_t^{\text{init}} = \underset{\{d_{\min} \dots d_{\max}\}}{\operatorname{argmin}} \sum_{\bar{x}} \left[ I - U(\bar{x}) f_d^{\text{init}}(\bar{x}, D_E(\bar{x})) + \sum_{\bar{y} \in N\bar{x}} \lambda(\bar{x}, \bar{y}) f_p(D_t(\bar{x}), D_t(\bar{y})) \right]$$

Goal: To preserve the discontinuities.  $\lambda$  is defined to encourage the disparity discontinuities to be consistent with intensity / color discontinuities.

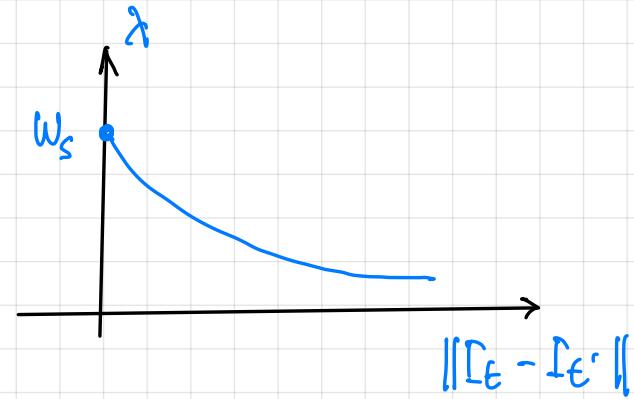
$$\lambda(\bar{x}, \bar{y}) = w_s \frac{U_\lambda(\bar{x})}{\|I_E(\bar{x}) - I_E(\bar{y})\| + \epsilon}$$

where:

$$U_\lambda(x) = \frac{|N_x|}{\sum_{\bar{y}' \in N\bar{x}} \frac{1}{\|I_E(\bar{x}) - I_E(\bar{y}')\| + \epsilon}}$$

$w_s$  = the smoothness strength

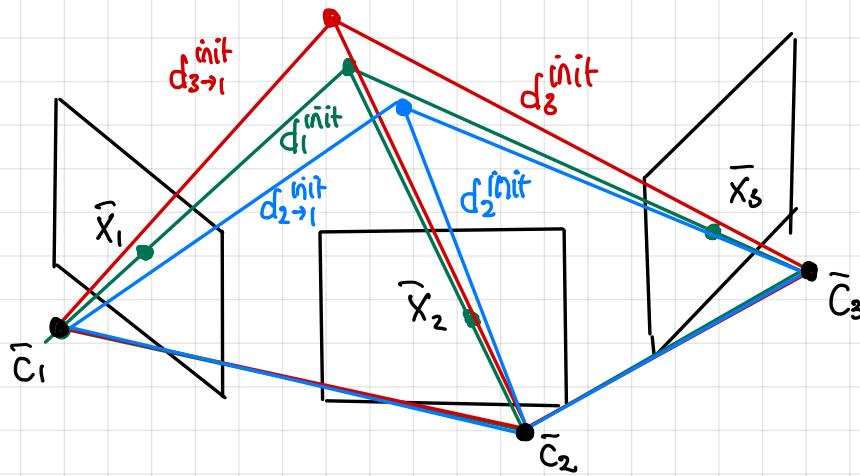
$$\lambda(\bar{x}, \bar{y}) = w_s \frac{|N_x|}{(\|I_E(\bar{x}) - I_E(\bar{y})\| + \epsilon) \sum_{\bar{y}' \in N\bar{x}} \frac{1}{\|I_E(\bar{x}) - I_E(\bar{y}')\| + \epsilon}}$$



1. The smaller the intensity/color difference ( $\|I_E(\bar{x}) - I_E(\bar{y})\|$ ) between two neighboring pixels, the larger  $\lambda$  is. A larger  $\lambda$  will encourage the smoothness constraint more.
2.  $U_\lambda$  is to normalize  $\|I_E(\bar{x}) - I_E(\bar{y})\|$  so that  $\lambda$  should be within a certain range depending on  $w_s$ .

## [8] Bundle Optimization

Problem :

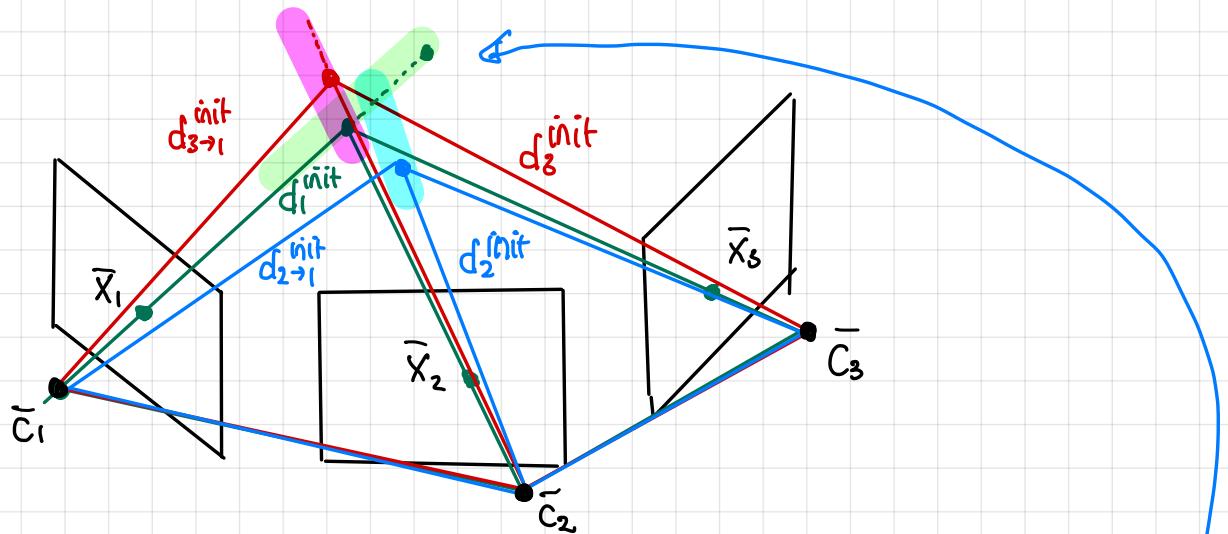


Ideally:  $d_1^{\text{init}} = d_{2 \rightarrow 1}^{\text{init}} = d_{3 \rightarrow 1}^{\text{init}}$

Reality:  $d_1^{\text{init}} \neq d_{2 \rightarrow 1}^{\text{init}} \neq d_{3 \rightarrow 1}^{\text{init}}$   $\rightarrow$  This is the cause of the flickering problem

To address the problem we need to make:

$d_1 \approx d_{2 \rightarrow 1} \approx d_{3 \rightarrow 1}$   $\rightarrow$  Enforce them to be as close as possible.



Basic idea :

To find other disparity value other than  $d_t^{\text{init}}$  that close the gap in  $d_1$  &  $d_{2 \rightarrow 1}$  &  $d_{3 \rightarrow 1}$

Candidates for  $d_t^{\text{init}}$  are:  $\{d_t^{\text{init}} - N\epsilon, \dots, d_t^{\text{init}}, d_t^{\text{init}} + \epsilon, \dots, d_t^{\text{init}} + N\epsilon\}$

$N$  &  $\epsilon$  are variables that need to be decided.

• Input:  $\{I_t, IP_t, D_t^{\text{init}}\}_{t=1}^N$

Output:  $D_t$

using the same functions as before

$$D_t = \underset{\{d_t^{\text{init}} + n\varepsilon\}_{n=-N}^{+N}}{\operatorname{argmin}}$$

$$\sum_{\bar{x}} \left[ \left( (1 - U(\bar{x}) f_d(\bar{x}, D_t(\bar{x})) \right) + \sum_{\bar{y} \in N\bar{x}} \lambda(\bar{x}, \bar{y}) f_p(D_t(\bar{x}), D_t(\bar{y})) \right]$$

candidates:  $\{d_{t-NE}^{\text{init}}, \dots, d_{t-\varepsilon}^{\text{init}}, d_t^{\text{init}}, d_{t+\varepsilon}^{\text{init}}, \dots, d_{t+NE}^{\text{init}}\}$

where:

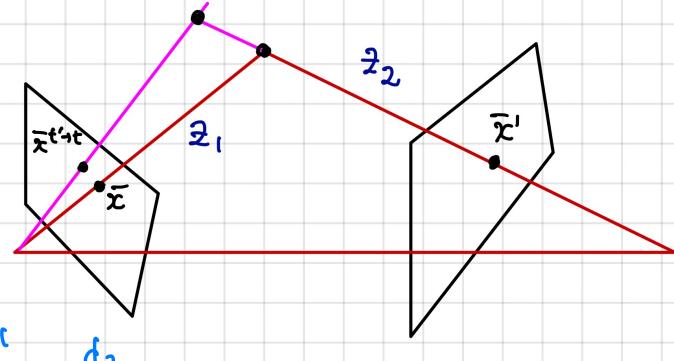
$$f_d(\bar{x}, d) = \sum_{t'}^N f_c(\bar{x}, d, I_t, I_{t'}) p_v(\bar{x}, d, D_{t'})$$

photo consistency      geometric coherence

$\Downarrow$   
the same as in the previous step.

Geometric Coherence:

Goal: to measure how consistent  $\bar{z}_1$  (or  $d_1$ ) and  $\bar{z}_2$  (or  $d_2$ ) geometrically.



$$p_v(\bar{x}, d, D_{t'}) = \exp \left( - \frac{\|\bar{x} - l_{t' \rightarrow t}(\bar{x}', D_{t'}(\bar{x}))\|^2}{2\sigma_d^2} \right)$$

where:  $\bar{x}' = l_{t' \rightarrow t}(\bar{x}, d)$

$$\bar{x}^{t' \rightarrow t} = l_{t' \rightarrow t}(\bar{x}', D_{t'}(\bar{x}))$$

$$\|\bar{x} - \bar{x}^{t' \rightarrow t}\| = \|\bar{x} - l_{t' \rightarrow t}(\bar{x}', D_{t'}(\bar{x}))\|$$

similar to  $\sigma_c$ : to control how tolerant we are with the error



;  $\sigma_d$  = standard deviation